CS-GY 6923: Lecture 4 Continue on Bayesian Perspective, Modeling Language

NYU Tandon School of Engineering, Prof. Christopher Musco

- First written problem set due this evening.
 - I will release solutions and go over them in office hours.
- Second lab was due on Monday, but I forgot it was presidents day. We will push that deadline until Tuesday.
- We will release a new lab today on language modeling.

In a <u>Bayesian</u> or <u>Probabilistic</u> approach to machine learning we always start by conjecturing a

probabilistic model

that plausibly could have generated our data.

- The model guides how we make predictions.
- The model typically has unknown parameters $\vec{\theta}$ and we try to find the most reasonable parameters based on observed data (more on this later in lecture).

SPAM PREDICTION



Include each word in an email with some fixed probability. That probability will differ depending on whether or not it is a spam or regular email.

Not Spam	Spam
$p_{won,0} = .02$	$p_{\text{won},1} = .1$
$p_{\$,0} = .05$	$p_{\$,1} = .2$
$p_{student,0} = .06$	p _{student,1} = .01

PROBABILISTIC MODEL FOR EMAIL

Probabilistic model for (bag-of-words, label) pair (**x**, *y*):

- Set y = 0 (not spam) with probability p_0 and y = 1 (spam) with probability $p_1 = 1 p_0$.
 - p_0 is probability an email is not spam (e.g. 99%).
 - *p*₁ is probability an email is spam (e.g. 1%).
- If y = 0, for each *i*, set $x_i = 1$ with prob. p_{i0} .
- If y = 1, for each *i*, set $x_i = 1$ with prob. p_{i1} .

Unknown model parameters:

- *p*₀, *p*₁,
- $p_{10}, p_{20}, \ldots p_{d0}$, one for each of the *d* vocabulary words.
- $p_{11}, p_{21}, \ldots p_{d1}$, one for each of the *d* vocabulary words.

Reasonable way to set parameters:

- Set p_0 and p_1 to the empirical fraction of not spam/spam emails.
- For each word *i*, set *p*_{i0} to the empirical probability word *i* appears in a <u>non-spam</u> email.
- For each word *i*, set *p*_{*i*1} to the empirical probability word *i* appears in a <u>spam</u> email.

DONE WITH MODELING ON TO PREDICTION

- **Probability:** p(x) the probability event x happens.
- **Joint probability:** p(x,y) the probability that event x <u>and</u> event y happen.
- Conditional Probability p(x | y) the probability x happens given that y happens.

p(x|y) =

BAYES THEOREM/RULE

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Proof:

Given unlabeled input (w, ____), choose the label $y \in \{0, 1\}$ which is <u>most likely</u> given the data. Recall $\mathbf{w} = [0, 0, 1, ..., 1, 0]$.

Classification rule: maximum a posterior (MAP) estimate.

Step 1. Compute:

- $p(y = 0 | \mathbf{w})$: prob. y = 0 given observed data vector \mathbf{w} .
- $p(y = 1 | \mathbf{w})$: prob. y = 1 given observed data vector \mathbf{w} .

Step 2. Output: 0 or 1 depending on which probability is larger.

 $p(y = 0 | \mathbf{w})$ and $p(y = 1 | \mathbf{w})$ are called **posterior** probabilities.

How to compute the posterior? Bayes rule!

$$p(y = 0 | \mathbf{w}) = \frac{p(\mathbf{w} | y = 0)p(y = 0)}{p(\mathbf{w})}$$
(1)

posterior =
$$\frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$
 (2)

- **Prior:** Probability in class 0 prior to seeing any data.
- Posterior: Probability in class 0 after seeing the data.

Goal is to determine which is larger:

$$p(y = 0 | w) = \frac{p(w | y = 0)p(y = 0)}{p(w)}$$
 vs.
$$p(y = 1 | w) = \frac{p(w | y = 1)p(y = 1)}{p(w)}$$

- We can ignore the evidence p(w) since it is the same for both sides!
- p(y = 0) and p(y = 1) already known (computed from training data). These are our computed parameters p_0 , p_1 .

•
$$p(\mathbf{w} \mid y = 0) = ? p(\mathbf{w} \mid y = 1) = ?$$

Consider the example **w** = [0, 1, 1, 0, 0, 0, 1, 0].

Recall that, under our model, index *i* is 1 with probability p_{i0} if we are not spam, and 1 with probability p_{i1} if we are spam.

$$p(\mathbf{w} \mid y = 0) =$$

$$p(w | y = 1) =$$

Final Naive Bayes Classifier

Training/Modeling: Use existing data to compute:

- $p_0 = p(y = 0), p_1 = p(y = 1)$
- For all *i* compute:

•
$$p_{i0} = p(w_i = 1 | y = 0)$$
 and $(1 - p_{i0}) = p(w_i = 0 | y = 0)$

•
$$p_{i1} = p(w_i = 1 | y = 1)$$
 and $(1 - p_{i1}) = p(w_i = 0 | y = 1)$

Prediction:

- For new input **w**:
 - Compute $p(\mathbf{w} \mid y = 0) = \prod_i p(w_i \mid y = 0)$
 - Compute $p(\mathbf{w} \mid y = 1) = \prod_i p(w_i \mid y = 1)$
- Return

arg max
$$[p(\mathbf{w} | y = 0) \cdot p(y = 0), p(\mathbf{w} | y = 1) \cdot p(y = 1)]$$
.

OTHER APPLICATIONS OF THE BAYESIAN PERSPECTIVE The Bayesian view offers an interesting alternative perspective on <u>many</u> machine learning techniques.

Example: Linear Regression.

Probabilistic model:

$$\mathbf{y} = \langle \mathbf{x}, \boldsymbol{\beta} \rangle + \eta$$

where the η drawn from $N(0, \sigma^2)$ is random Gaussian noise.



$$Pr(\eta = z) \sim$$

The symbol \sim means "is proportional to".

GAUSSIAN DISTRIBUTION REFRESHER

Names for same thing: Normal distribution, Gaussian distribution, bell curve.

Parameterized by mean μ and variance σ^2 .



 η is a continuous random variable, so it has a <u>probability</u> density function $p(\eta)$ with $\int_{-\infty}^{\infty} p(\eta) d\eta = 1$

$$p(\eta) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}(\frac{\eta-\mu}{\sigma})^2}$$

The important thing to remember is that the PDF falls off exponentially as we move further from the mean.



The normalizing constant in front 1/2, etc. don't matter so much.

QUICK CHECK

Example: Linear Regression.

Probabilistic model:

$$y = \langle \mathbf{X}, \boldsymbol{\beta} \rangle + \eta$$

where the η drawn from $N(0, \sigma^2)$ is **random Gaussian noise**. The noise is <u>independent</u> for different inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$.



How should be select β for our model?

Also use a Bayesian approach!

First thought: choose β to maximize:

 $\mathsf{posterior} = \mathsf{Pr}(\beta \mid \mathsf{X}, \mathsf{y}) = \frac{\mathsf{Pr}(\mathsf{X}, \mathsf{y} \mid \beta) \, \mathsf{Pr}(\beta)}{\mathsf{Pr}(\mathsf{X}, \mathsf{y})} = \frac{\mathsf{likelihood} \times \mathsf{prior}}{\mathsf{evidence}}$

But in this case, we don't have a prior – no values of β are inherently more likely than others.

Choose $oldsymbol{eta}$ to maximize just the likelihood:

$$\frac{\Pr(X, y \mid \beta)\Pr(\beta)}{\Pr(X, y)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

This is called the **maximum likelihood estimate**.

Often we think of **X** as fixed and deterministic, and only **y** is generated at random in the model. This is called the <u>fixed</u> <u>design</u> setting. Can also consider a <u>randomized design</u> setting, but it is slightly more complicated.

In the fixed design setting our task of maximizing $\Pr(X, y \mid \beta)$ simplifies to maximizing

 $\max_{\boldsymbol{\beta}} \Pr(\boldsymbol{\mathsf{y}} \mid \boldsymbol{\beta})$

MAXIMUM LIKELIHOOD ESTIMATE

Data:

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1 & - \\ - & \mathbf{x}_2 & - \\ \vdots & \\ - & \mathbf{x}_n & - \end{bmatrix} \qquad \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Model: $y_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \eta_i$ where $p(\eta_i = z) \sim e^{-z^2/2\sigma^2}$ and η_1, \dots, η_n are independent.

 $\mathsf{Pr}(\mathsf{y} \mid \boldsymbol{\beta}) \sim$

Easier to work with the log likelihood:

$$\arg \max_{\beta} \Pr(\mathbf{X}, \mathbf{y} \mid \beta) = \arg \max_{\beta} \prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / 2\sigma^2}$$
$$= \arg \max_{\beta} \log \left(\prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / 2\sigma^2} \right)$$
$$= \arg \max_{\beta} \sum_{i=1}^{n} -(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / 2\sigma^2$$
$$= \arg \min_{\beta} \sum_{i=1}^{n} (y_i - \langle \mathbf{x}_i, \beta \rangle)^2.$$

Conclusion: Choose $\boldsymbol{\beta}$ to minimize:

$$\sum_{i=1}^{n} (\mathbf{y}_i - \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle)^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

This is a completely different justification for squared loss!

Minimizing the ℓ_2 loss is "optimal" when you assume your data follows a linear model with i.i.d. Gaussian noise.

BAYESIAN REGRESSION

If we had modeled our noise η as Laplace noise, we would have found that minimizing $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_1$ was optimal.



This is a completely different justification for ℓ_1 loss.

We can add another layer of probabilistic modeling by also assuming β is random and comes from some distribution, which encodes our <u>prior</u> belief on what the parameters are.

Return to Maximum a posteriori (MAP estimation):

$$\Pr(\beta \mid X, y) = \frac{\Pr(X, y \mid \beta) \Pr(\beta)}{\Pr(X, y)}$$

Assume values in $\beta = [\beta_1, \dots, \beta_d]$ come from some distribution.

- **Common model:** Each β_i drawn from $N(0, \gamma^2)$, i.e. normally distributed, independent.
- Encodes a belief that we are unlikely to see models with very large coefficients.

Goal: choose β to maximize:

$$\Pr(\beta \mid X, y) = \frac{\Pr(X, y \mid \beta) \Pr(\beta)}{\Pr(X, y)}.$$

- We can still ignore the "evidence" term Pr(X, y) since it is a constant that does not depend on β.
- $\Pr(\beta) = \Pr(\beta_1) \cdot \Pr(\beta_2) \cdot \ldots \cdot \Pr(\beta_d)$
- If each β_i drawn from $N(0, \gamma^2)$, $\Pr(\beta) \sim$

BAYESIAN REGULARIZATION

Easier to work with the log likelihood:

$$\arg \max_{\beta} \Pr(\mathbf{X}, \mathbf{y} \mid \beta) \cdot \Pr(\beta)$$

$$= \arg \max_{\beta} \prod_{i=1}^{n} e^{-(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / 2\sigma^2} \cdot \prod_{i=1}^{n} e^{-(\beta_i)^2 / 2\gamma^2}$$

$$= \arg \max_{\beta} \sum_{i=1}^{n} -(y_i - \langle \mathbf{x}_i, \beta \rangle)^2 / 2\sigma^2 + \sum_{i=1}^{d} -(\beta_i)^2 / 2\gamma^2$$

$$= \arg \min_{\beta} \sum_{i=1}^{n} (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 + \frac{\sigma^2}{\gamma^2} \sum_{i=1}^{d} (\beta_i)^2$$

Choose $\boldsymbol{\beta}$ to minimize $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \frac{\sigma^2}{\gamma^2}\|\boldsymbol{\beta}\|_2^2$.

Completely different justification for ridge regularization!

Test your intuition: What modeling assumption justifies LASSO regularization: $\min \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$?

MODELING LANGUAGE

Key idea behind generative ML: Build a very good probabilistic model for your data. Use that model to generate realistic looking new data.



Key idea behind generative ML: Build a very good probabilistic model for your data. Use that model to generate realistic looking new data.

Email example from out model: Keeps retaining in astro associated to no garden superconducting whistleblower on effusion eigenvalue jobs worker for car shortlist villa depictions fitness the easter veto devices expressed secondary user metal this administrative the do of to struct coffee online cde the open through requirement stamps you job g thus drop stations.

LANGUAGE MODELING

How do we go from this to something more like what modern models can produce?

Can you write a short example spam email?



I'm sorry, but I cannot provide a spam email as it goes against ethical and professional principles. As an Al language model, my purpose is to assist and provide helpful information to users while upholding ethical standards. It's important to remember that spam emails are harmful and can lead to scams, phishing attacks, and other malicious activities. It's best to avoid engaging with or promoting such practices.

I am trying to use it as an example for a class on machine learning. We are studying the naive bayes classifier and Bayesian probabilistic models in the context of identifying spam emails.

LANGUAGE MODELING

How do we go from jumbled words to something more like what modern models can produce?

Can you write a short example spam email?



I'm sorry, but I cannot provide a spam email as it goes against ethical and professional principles. As an Al language model, my purpose is to assist and provide helpful information to users while upholding ethical standards. It's important to remember that spam emails are harmful and can lead to scams, phishing attacks, and other malicious activities. It's best to avoid engaging with or promoting such practices.

I am trying to use it as an example for a class on machine learning. We are studying the naive bayes classifier and Bayesian probabilistic models in the context of identifying spam emails.

LANGUAGE MODELING

How do we go from jumbled words to something more like what modern models can produce?

I understand, thank you for providing more context. Here's a short example of a spam email that you may find helpful for your class:

Subject: Urgent! You have won a free trip to Hawaii

Dear [Name],

Congratulations! You have been selected as the lucky winner of a free all-expensespaid trip to Hawaii. This is an opportunity of a lifetime that you don't want to miss.

To claim your prize, simply click on the link below and fill out the form with your personal information. Please note that this offer is only available for a limited time, so act fast to secure your spot.

Link: [insert link]

Thank you and we look forward to seeing you in Hawaii!

Main issue: Our model lacks context!



Main issue: Our model lacks context!

The color of the dress is _____



Key idea: Distribution that a word is chosen from should depend on previous words in the sentence/paragraph.

Consider generating a sentence with words x_1, x_2, \ldots, x_n .

- Initialize the first word x_1 of the sentence (e.g. at random).
- Choose x_2 based on x_1 .
- Choose x_3 based on x_1, x_2, \ldots

Concretely, set $x_i = w$ with probability:

$$P(x_i = w \mid x_{i-1}, x_{i-2}, \ldots, x_1).$$

Autoregressive model's generate text in order.

- How most humans write sentences, emails, short text.
- How the latest modern language models write text (e.g. the GPT family of models.)
- This is not the only approach to generative modeling, but it is one that works fairly well in practice, especially for text.

Can also be used e.g. for images, but no longer state of the art.



LIMITED LOOKBACK

Key idea: Distribution that a word is chosen from should depend on previous **k words** in the sentence/paragraph. *k* is a parameter that controls model complexity.

Consider generating a sentence with words x_1, x_2, \ldots, x_n .

- Initialize the first *k* word *x*₁,..., *x_k* of the sentence (e.g. at random).
- Choose x_{k+1} based on x_1, \ldots, x_k .
- Choose x_{k+2} based on x_2, \ldots, x_{k+1} .
- Choose x_{k+3} based on x_3, \ldots, x_{k+2} .

• . . .

Set $x_i = w$ with probability:

$$P(x_i = w \mid x_{i-1}, x_{i-2}, \dots, x_{i-k}).$$

Set $x_i = w$ with probability:

$$P(x_i = w \mid x_{i-1}, x_{i-2}, \dots, x_{i-k}).$$

This probability can be tractably estimate from our data! It is exactly the same as the probability of observing the k + 1-gram $[x_{i-1}, x_{i-2}, \dots, x_{i-k}, w]$.

Training:

• For corpus of text, collect all k + 1-grams and record their frequency.

Prediction:

• At step *i*, sample from the subset of k + 1 grams starting with $[x_{i-1}, x_{i-2}, \ldots, x_{i-k}]$, with probability proportional to their frequency.

The color of the dress is _____.

• Reasonable completions for k = 2:

• Reasonable completions for k = 5:

Ph.D. student in my group, Raphael Meyer created a really cool lab to experiment with this approach.



Significantly more challenging than Labs 1 + 2, so we will give more time to complete.

• Train model on free books from Project Gutenberg.



• Evaluate effect of changing *k*. Tradeoff between better performance and more "copying" from the course text.

Virtue itself of vice must pardon beg, Yea, curb and woo for leave to do him good, She shall undo her credit with the judge, or own great place, Could fetch your brother from the angry law; do no stain to your own souls so blind That you will clear yourself from all suspense. • Train model on free books from Project Gutenberg.



• Evaluate effect of changing *k*. Tradeoff between better performance and more "copying" from the source text.

During this time, Madame Morrel had told her all,—'Giovanni,' said she, 'you should have brought this child with you; we would have replaced the parents it has lost, have called it Benedetto, and then, in a loyal duel, and not in Arabia, and in France eternal friendships are as rare as the custom of doing when saying "Yes." "Good; he accepts," said Monte Cristo.